



Dent, Ian and Craig, Tony and Aickelin, Uwe and Rodden, Tom (2014) Variability of behaviour in electricity load profile clustering: who does things at the same time each day? In: Advances in data mining: applications and theoretical aspects: 14th Industrial Conference, ICDM 2014, St. Petersburg, Russia, July 16-20, 2014: proceedings. Lecture notes in computer science (8557). Springer International Publishing, Cham, pp. 70-84. ISBN 9783319089768 (electronic bk.); 9783319089751 (print)

Access from the University of Nottingham repository:

http://eprints.nottingham.ac.uk/3347/1/ICDM_2014.pdf

Copyright and reuse:

The Nottingham ePrints service makes this work by researchers of the University of Nottingham available open access under the following conditions.

- Copyright and all moral rights to the version of the paper presented here belong to the individual author(s) and/or other copyright owners.
- To the extent reasonable and practicable the material made available in Nottingham ePrints has been checked for eligibility before being made available.
- Copies of full items can be used for personal research or study, educational, or not-for-profit purposes without prior permission or charge provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way.
- Quotations or similar reproductions must be sufficiently acknowledged.

Please see our full end user licence at:

http://eprints.nottingham.ac.uk/end_user_agreement.pdf

A note on versions:

The version presented here may differ from the published version or from the version of record. If you wish to cite this item you are advised to consult the publisher's version. Please see the repository url above for details on accessing the published version and note that access may require a subscription.

For more information, please contact eprints@nottingham.ac.uk

Variability of Behaviour in Electricity Load Profile Clustering; Who Does Things at the Same Time Each Day?

Ian Dent¹, Tony Craig², Uwe Aickelin¹, and Tom Rodden¹

¹ School of Computer Science, University of Nottingham, Nottingham NG8 1BB, UK,
psxid@nottingham.ac.uk,

WWW home page: <http://ima.ac.uk/dent>

² The James Hutton Institute, Aberdeen, UK

Abstract. UK electricity market changes provide opportunities to alter households' electricity usage patterns for the benefit of the overall electricity network. Work on clustering similar households has concentrated on daily load profiles and the variability in regular household behaviours has not been considered. Those households with most variability in regular activities may be the most receptive to incentives to change timing. Whether using the variability of regular behaviour allows the creation of more consistent groupings of households is investigated and compared with daily load profile clustering. 204 UK households are analysed to find repeating patterns (motifs). Variability in the time of the motif is used as the basis for clustering households. Different clustering algorithms are assessed by the consistency of the results.

Findings show that variability of behaviour, using motifs, provides more consistent groupings of households across different clustering algorithms and allows for more efficient targeting of behaviour change interventions.

1 Background and Motivation

The electricity market in the UK is undergoing dramatic changes. Legal, social and political drivers for a more carbon efficient electricity network, along with the dramatically increased flow of data from households through the deployment of smart meters, requires a transformation of existing practices. In particular, the change of the frequency of sampling of electricity usage, by using smart meters, alters the level of understanding of households' behaviour that is possible [1].

One approach to address the pressures on the electricity network is the application of Demand Side Management (DSM) techniques to achieve changes in consumer behaviour. DSM is defined as “systematic utility and government activities designed to change the amount and/or timing of the customer’s use of electricity” for the collective benefit of society, the utility company, and its customers [2]. The peak time for electricity usage in the UK is during the early evening and the successful application of techniques to reduce, or move, the peak usage would improve the overall efficiency of the electricity network.

To allow selection of appropriate DSM interventions, a good understanding of the existing behaviour of households is needed. Firstly, knowledge is needed on an individual household that can be deduced from house-wide electricity metering. Secondly, a method is required to group large numbers of households into a manageable number of archetypal groups where the members display similar characteristics. This approach allows for cost effective targeting of the most appropriate subset of customers whilst allowing the company management to deal with a manageable number of archetypes [3].

There is an extensive body of work on clustering households which includes comparing or combining timed meter readings to create additional attributes that contribute to the quality of the clustering [4]. However, little work has focused on how the daily activity patterns of the household vary from day to day and how this can be used for clustering. For instance, some households will be creatures of habit and will eat their evening meal at almost the same time each evening, whilst others have a much more variable activity pattern and will eat at different times. Ellegård and Palm [5] have investigated the variability of behaviour using diaries and interviews but have not used analysis of meter data.

Clustering households using their degree of variability in behaviour, as shown by electricity consumption, provides a way of identifying the subset of electricity users who may be most receptive to an intervention to influence their activity patterns. The intervention may be to reward households for NOT changing their current pattern of usage if it is already as desired by the utility company.

This paper addresses the question of whether making use of the variability of behaviour (as shown by the electricity meter data) provides “better” groupings of households for the purpose of DSM than those provided by using daily load profiles. The judgement of “better” is measured by implementing a number of different clustering techniques and measuring the degree of overlap between the clusters found. A consistent set of clusters across the different clustering algorithms implies a better, and more useful, approach to generating the clusters.

The investigation of household electricity load profiles is an important area of research given the centrality of such patterns in directly addressing the needs of the electricity industry, both now and in the future. This work extends existing load profile work by taking electricity meter data streams and developing new ways of representing the household that can be used as the basis for clustering using existing data mining techniques. The identification of repeating motifs and the investigation of how the timing of the motifs varies from day to day, as a key behavioural trait of the household, is a novel area of research. An improvement in creating useful archetypes can have major financial and environmental benefits.

2 Methods and Technical Solutions

2.1 Load Profiling

There has been extensive research on determining daily load profiles to represent a household’s electricity usage [6]. In many cases, (e.g., [7]), the daily load profiles

are used as the basis for clustering “similar” households together to develop a small set of archetypal profiles which can be used for targeting of behaviour change interventions. Previous work has used different clustering techniques with the majority of the published literature using hierarchical clustering.

The common approach is to define a subset of the data (e.g. by season and/or by day of the week) and then to create average daily profiles for a household from the electricity meter data. The shapes of these daily profiles are then clustered to group similar shapes together. A representative profile is defined (e.g. by averaging all the members of the cluster) to produce a archetypal daily load profile for that cluster of households.

Previous work has not investigated how households may exhibit different behaviour from day to day and how these differences may be used as a distinguishing feature of the household and a basis for clustering.

2.2 Motifs

The electricity meter data reading stream from a household can be plotted as a graph of usage against time and regular activities appear as similar shaped patterns. Short patterns that repeat are defined as “motifs” and detection of these motifs, and their timing, can inform understanding of household behaviour.

This work uses the SAX (Symbolic Aggregate approXimation) technique which allows symbolic representation of time series data [8, 9]. Other motif finding algorithms could also be incorporated into the proposed approach to identify the flexibility of behaviour (e.g. [10]). To assess variability within a household, it is necessary to detect the repeating motifs that are assumed to signify particular activities (e.g., cooking the evening meal). These are generally of a similar shape on different days but show some differences due to noise caused by other activities within the household (e.g., a fridge automatically running). The SAX approach of symbolising the real valued meter readings is useful as it allows for approximate matching (as various ranges of readings map to a single symbol).

Lines et al [11] applies motif finding to UK data to detect the use of particular appliances, drawn from a set of known appliances. This contrasts with the focus in this paper which is to find interesting, repeating patterns of behaviour without the need to define the activity that the motif represents. Appliances that can be consistently and accurately detected can be used with the approach detailed here by extending the analysis of the timing of repeating motifs to the analysis of variability of timing of appliance usage.

2.3 North East Scotland Electricity Monitoring Project (NESEMP)

This study makes use of data collected as part of the ongoing NESEMP which is examining the relationship between different types of energy feedback and psycho-social measures including individual environmental attitudes, household characteristics, and everyday behaviours. As part of this ongoing project, several hundred households are being monitored and the electricity usage is recorded every five minutes using CurrentCost monitors [12].

After removing data for households with insufficient readings, the data is loaded into a MySQL database and the readings are aligned with exact 5 minute boundaries (e.g. 1pm, 1.05pm, etc.) by interpolation between the actual readings. This is achieved by calculating the reading at an exact 5 minute point (e.g. 1.05pm) by considering the actual readings before and after that time and by calculating the reading such that the total usage over a longer period is the same whether the interpolated readings or the original actual readings are used [13]. This results in a set of 288 readings (one for every 5 minute period in the day) for each of the households in the database.

Each day of sampling is labelled in a number of ways such as “working day” or “summer” to aid selection of particular subsets of data.

2.4 Detecting Motifs

To find motifs within the data, each period of interest within the day (e.g., the peak period) for each household is examined by taking a moving window over the period. The subset of the meter readings within the moving window is then converted into a string and stored. Next, the window moves on by one time period (5 minutes) and the conversion into a string is repeated. Using an alphabet size of 5 and a motif size of 6 (i.e., 30 minutes), analysing the 4pm to 8pm period provides a total of 49 x 5 minute readings for each day. As the interest is in changes in usage rather than absolute usage, these readings are compared with adjacent readings in time to produce 48 values (one for each 5 minute period) representing the change in usage since the last 5 minute reading. This results in 42 motifs stored for each day for each household (one for each possible 30 minute period within the peak time). Fig. 1 shows an example of how the symbolised motifs are built up. The top graph shows the 5 minute readings for the 4 hour peak period. A sliding window of 6 readings (30 minutes) is taken across the peak period with the first 2 and the last window shown. Each window is normalised within the values in the window and then translated into the symbolised representations as shown at the bottom of the diagram.

The analysis uses an alphabet of 5 symbols (i.e., the letters “a” to “e”) to represent the motifs. 5 is selected as a reasonable compromise between having too few symbols, and thus not detecting changes in electricity consumption, and having too many and thus generating too many patterns that do not repeat. The symbolisation translates readings within a particular range into a given letter and thus similar, although not identical, readings are translated into the same letter. The resulting motifs for 2 windows may be identical whereas the original readings may only be approximately similar.

The motif size selected is 6 corresponding to a 30 minute (i.e., 6 x 5 minutes) period. This figure was selected as the UK electricity settlement market uses a 30 minute period [14] and 30 minutes is also a reasonable period that will allow time for activities such as showering.

The motifs are built from the graph shape without regard to absolute value of the data. A possible effect of this is to find motifs within what is the general

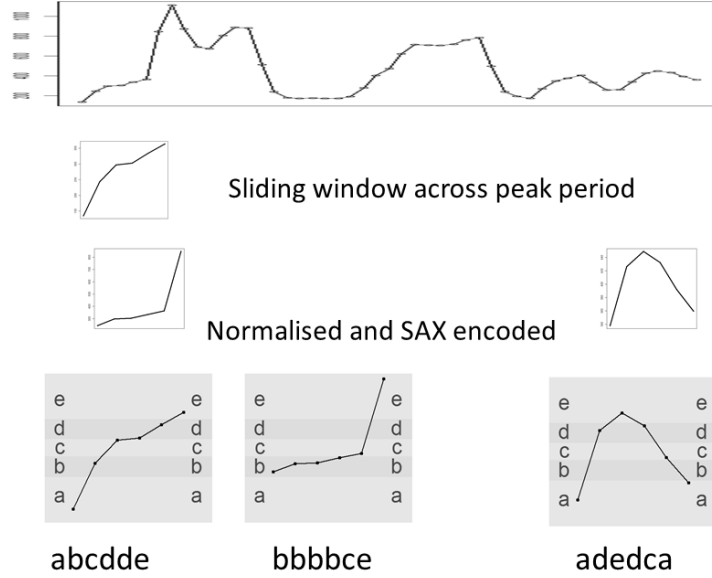


Fig. 1. Example of symbolisation (alphabet of 5, motif length of 6)

noise associated with the meter readings. This is avoided by ignoring any motifs within a window which have a range of less than 100W.

As the motifs are created by shifting a moving window over the stream of data, overlapping periods are considered and periods with no activity except for one change in meter reading will lead to a series of motifs that are similar. For example, a long period of no activity except for a jump of +200W will lead to motifs being found such as ccccca, ccccac, cccacc, etc. As only one of these is interesting for further analysis, the others are excluded.

The top motif (the one that occurs most often within a household) is further examined for the times when the motif occurs on each day. The number of times the motif occurs, and the standard deviation of the time of occurrence, are calculated for each household. Similarly, the second and third most common motifs within a household are identified and the variability in timing calculated.

Other useful measures relating to the motifs found within a household are also calculated including the number of different motifs (occurring at least twice) and the number of different motifs occurring on at least 30% of the days sampled for the household. The 30% figure is selected as a reasonable number to ensure only regularly repeating patterns are considered.

The attributes calculated for each household and used as input to the clustering algorithms are:

1. Number of occurrences of the motif occurring most frequently during the peak period.
2. Variability in timing of the occurrence of the most frequent motif within the household. This is represented by the standard deviation of the timing (measured in minutes) around the mean start time.
3. Number of occurrences of the second most frequent motif.
4. Variability in timing of the second most frequent motif.
5. Number of occurrences of the third most frequent motif.
6. Variability in timing of the third most frequent motif.
7. Total number of motifs for the household that occur at least twice.
8. Total number of different motifs that occur on at least 30% of days.

2.5 Clustering Algorithms

Various clustering techniques are selected for evaluation of the different approaches to analysing the data. Note that, whilst possibly a useful additional benefit, this work does not focus on selecting the “best” clustering algorithm but uses a selection of algorithms to assess the benefits or otherwise of making use of the motif variability information.

Based on the review by Chicco [6] the following clustering algorithms are selected as the most commonly used in previous work:

1. Kmeans is a well known algorithm that occurs in a number of examples of previous load profiling work. The algorithm requires a number of clusters (k) and works by randomly selecting an initial k locations for the centres of the clusters. Each data point is then assigned to one of the clusters by selecting the centre nearest to that data point. Once all the data points are assigned, each collection of points is considered, the new centre of the allocated points is calculated and the centre for that cluster is reassigned. The points are then reallocated to their new nearest centre and the algorithm continues until no changes are made to the allocations of points for an iteration [15].
2. Fuzzy c means. This provides an extension of the kmeans algorithm allowing partial membership to more than one cluster. The algorithm provides additional output showing the degree of membership that each household has of each of the derived clusters [16]. For this analysis, each household is assigned to the cluster for which they have the highest degree of membership.
3. Self Organising Maps. The Self Organising Map (SOM) is a neural network algorithm that can be used to map a high dimension set of data into a lower dimension representation. In this paper, the mapping is to a 2 dimensional set of representations which are arranged in a hexagonal map. Each sample (e.g., the average load profile for a given household) is assigned to a position in the map depending on the closeness of the sample to the existing nodes assigned to each position in the map (using a Euclidean measure of distance). Initially the nodes are assigned at random but, over time, the map produces an arrangement where similar samples are placed closely together and dissimilar samples are placed far apart [17].

4. Hierarchical clustering. Most of the published load profiling work has used hierarchical clustering and this approach has the benefit of providing easily understood rules for cluster membership. The algorithm uses a dissimilarity matrix for the households and, starting initially with each household in its own cluster, proceeds by joining clusters which are most similar. The hierarchy is cut at a point to provide the desired number of clusters [18]. The Euclidean distance is used when creating the dissimilarity matrix and the Ward agglomeration method [19] is used for combining clusters. The Ward method minimises the sum of squares of possible clusters when selecting households to combine. Other agglomeration techniques tend to create a few small clusters containing extreme valued households plus one large cluster containing the remainder of the households.
5. Random Forests [20] is used to create a dissimilarity matrix which is used with Partitioning Around Medoids (pam) to form clusters. This is implemented using the R package randomForest [21].

A common issue is the appropriate setting for the number of clusters. To match common practice within the electricity industry, 8 clusters are selected. The UK electricity industry has worked with 8 load profiles since the 1990s [22]. Figueiredo et al [23] report that the Portuguese electricity utility aim for a number of clusters between 6 and 9.

2.6 Cluster Validity Measures

To assess the benefits of a particular cluster solution an appropriate cluster validity index needs to be used. Many have been considered in the literature with the Mean Index Adequacy (MIA) and the Cluster Dispersion Indicator (CDI) [24] used in most of the published load profiling work. Lower values for the CDI and MIA measure denote “better” solutions.

The data to be clustered consists of M records numbered as $m = 1, \dots, M$. Each record has H attributes numbered as $h = 1, \dots, H$. The h th attribute for the i th record is designated as $m_i(h)$.

The data is clustered into K clusters (numbered as $k = 1, \dots, K$). Each cluster has R_k members where $r_{(k)}$ is the r th record assigned to cluster k and $C_{(k)}$ is the calculated centre of the cluster k .

The distance (d) between 2 records is defined as:

$$d(m_i, m_j) = \sqrt{\frac{1}{H} \sum_{h=1}^H (m_i(h) - m_j(h))^2} \quad (1)$$

where $m_i(h)$ and $m_j(h)$ are the h th attributes for two records, m_i and m_j .

The “within set distance” $\hat{d}(S)$ of the members of a set, S with N members (s_j where $j = 1, \dots, N$) is defined as:

$$\hat{d}(S) = \sqrt{\frac{1}{2N} \sum_{n=1}^N \sum_{p=1}^N d^2(s_n, s_p)} \quad (2)$$

The MIA gives a value which relies on the amount by which each cluster is compact - i.e., if the members in the cluster are close together the MIA is low.

$$MIA = \sqrt{\frac{1}{K} \sum_{k=1}^K \sum_r d^2(r_{(k)}, C_{(k)})} \quad (3)$$

The CDI depends on the distance between the members of the same cluster (as for the MIA) but also incorporates information on the distances between the representative load diagrams (i.e., the centroids) for each cluster. This therefore measures both the compactness of the clusters and the amount by which each cluster differs from the others.

$$CDI = \frac{1}{\hat{d}(C)} \sqrt{\frac{1}{K} \sum_{k=1}^K \hat{d}^2(R_k)} \quad (4)$$

where C is the set of cluster centres and R_k is the k th cluster members set.

2.7 Processing

UK specific data is used to generate average daily load profiles for each household which are clustered to provide a baseline for comparison. Selected clustering algorithms are applied to the data and validity indexes are used to produce a measure of the quality of the partitions found.

Next, the novel approach of identifying motifs within the data, and measuring the variability in timing of the motifs, is used to generate a new set of derived data using the same UK dataset. The same clustering algorithms and validity indexes are then applied to this dataset. In addition, the results are compared with the baseline obtained from the average daily load profiles in the first step.

2.8 Assessing the Results

To assess the consistency of clustering solutions, the different arrangements of households into clusters are compared. The consistency of the clusters obtained from the different clustering algorithms is used as a measure of the quality of the results with more consistency between the results suggesting a more useful method of identifying the clusters.

Measuring consistency across the clustering results using the different sets of data (load profiles and motifs) may be criticised as not necessarily providing a true measure of quality as clustering results may be consistent but not necessarily represent useful, “true” clusters within the data.

The Rand index compares the different pairs of samples (i.e., each possible pair of households) and assesses the number in which each pair are in the same partition in the 2 different clustering solutions, the number where each member of the pair are in different partitions in both solutions, and the case where the members are in the same partition in one solution but a different partition in

the other solution. The corrected Rand index [25] builds on the original work but adjusts the calculated value for the expected matching that would occur in a random arrangement. The corrected Rand index ranges from -1 to 1 with a higher value signifying better agreement between the partitions and hence a better solution.

3 Empirical Evaluation

3.1 Data Selection

A subset of the data is extracted for the peak period of 4pm to 8pm and for working days from Spring (March, April and May) 2011. Working days are weekdays excluding Scottish public holidays. Not all households have a full set of meter readings and those with less than 4 days of valid readings are excluded. The dataset has around 440,000 individual meter readings from 204 households.

The activities of interest within a household are related to switching appliances on or off (e.g., the use of electrical appliances in cooking) and it is the changes in the readings, rather than the absolute readings, that are of most interest and are used as the basis for analysis when using motifs.

3.2 Clustering Using the Load Profile Data

The data for the evening peak period (4pm to 7.55pm) are averaged to create a representative load profile for each household. For example, all the readings for 4pm for the household are averaged to create a representative reading for 4pm, similarly for 4:05pm, etc. The 204 representative profiles, each with 48 attributes (one for each time point), are then normalised within the 0 to 1 range and used with a variety of clustering algorithms.

3.3 Non-Motif Variability Clustering

Various different measures of variability of behaviour within the household can be defined without the use of motifs (e.g., [26]) and two methods are considered.

One approach is to consider the time at which the maximum usage occurred on each day during the period of analysis. These times are then used to calculate the standard deviation of the time around the mean for each household.

A second approach is to consider the total usage during the peak period on each day during Spring 2011. The standard deviation of the total per day around the mean total per day also provides a measure of variability of behaviour. Each of the 2 measures are calculated and used as the basis of simple clustering using kmeans ($k = 8$). The households in each of the clusters are shown in Fig. 2.

There is little correspondence between the cluster assignments for the 2 methods. The corrected Rand index of 0.01 shows no correspondence beyond that expected by chance. Furthermore, there is little correspondence with the clusters obtained from the motif variability approach (detailed below). A Spearman's rho

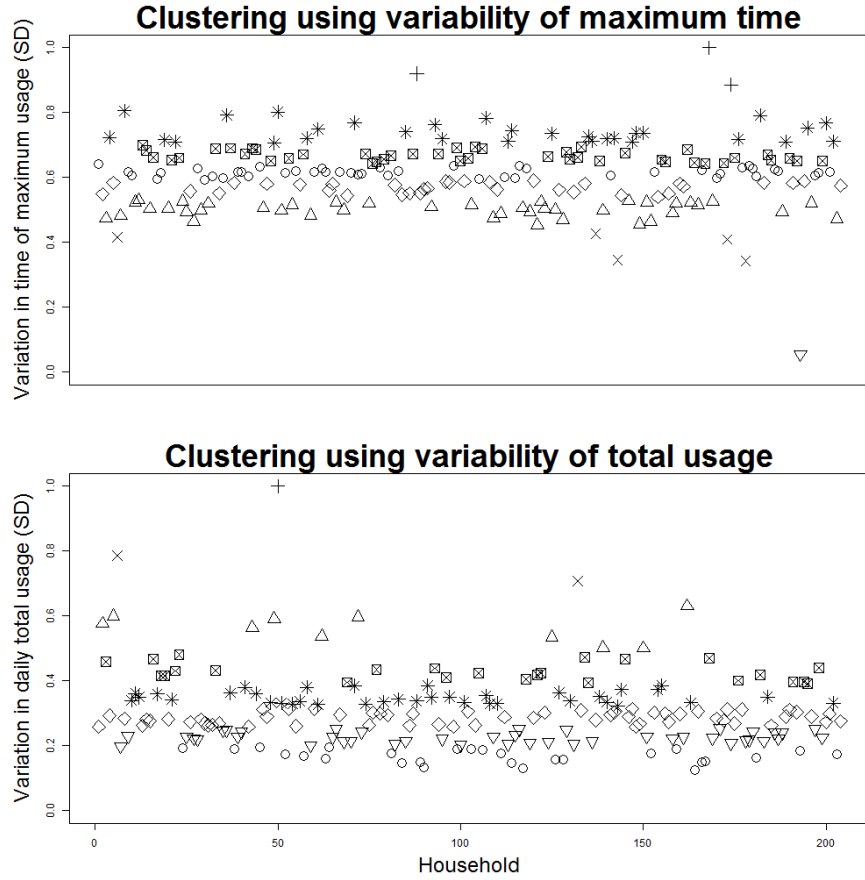


Fig. 2. Results for alternative variability measures

value of 0.23 shows there is little correlation between the two different variability measures.

It is therefore concluded that neither of the non-motif measures give a useful, consistent measure of the variability of each household.

3.4 Clustering Using the Motif Data

This paper finds the motifs in the stream of meter data and then examines how the times of these repeating patterns vary from day to day within a household. Furthermore, the number of times a pattern repeats within a household is also used as an indication of the variability of behaviour of that household.

The motifs in the data are discovered and the attributes detailed in Section 2.4 are generated. The same clustering algorithms as used for the load profile clustering are then applied to produce 8 archetypal clusters.

3.5 Results

Various measures that represent the variability of behaviour can be constructed and this paper considers the variability in time of maximum usage and the variability in total usage. However, as each measure is intended to represent the same thing (i.e., the variability of behaviour), the fact that there is little correlation between the measures, or the membership of the clusters generated using the measures, means that they provide a poor representation of the characteristic.

Comparing the load profile results with the motif variability results, Table 1 shows, for each of the clustering algorithms used and for each set of data, the sizes of the partitions in the solution and the values for the MIA and CDI cluster validity indexes (lower is better).

Table 1. Clustering Results and Validity indexes

	Load Profiles			Motifs		
	Cluster sizes	MIA	CDI	Cluster sizes	MIA	CDI
Kmeans	10,16,19,20,20,27,44,48	0.593	1.34	2,5,7,26,29,37,41,57	0.445	0.641
Fuzzy	14,17,20,23,23,23,35,49	0.679	2.14	12,15,19,26,28,30,34,40	0.551	2.084
SOM	13,15,16,20,28,31,36,45	0.595	1.337	2,5,24,25,28,29,40,51	0.451	0.733
Hier	9,10,13,20,22,37,43,50	0.61	1.386	2,3,5,26,31,34,40,63	0.46	0.64
RF	14,18,19,28,29,29,30,37	0.794	1.131	18,18,19,21,25,32,35,36	0.628	1.34

The MIA and CDI values show that the kmeans and SOM techniques produce similar quality solutions using the load profiles. The hierarchical algorithm is less good with the Fuzzy Cmeans algorithm being significantly poorer. The random forest and pam combination provides a good result for the CDI measure but scores poorly on the compactness of the clusters (as measured by MIA).

When using the motif variability data, the kmeans, SOM and hierarchical algorithms produce similar quality results with the Fuzzy Cmeans algorithm again producing poorer results. The random forest and pam combination provides middling results.

The MIA and CDI validity index calculations are not comparable between datasets due to the different number of attributes used.

Table 2 gives information on the consistency of the cluster partitions as the clustering algorithm changes. The results for the Rand index show that the values are consistently closer to 1 in the case of the clusters built using motif variation. The mean values for the Rand index (after omission of the values on the diagonal) are 0.4549 for the load profiles and 0.5183 for the motif variability approach. This shows a more consistent set of partitions are created when using the motif variability than the partitions created using the load profile information.

Table 2. Modified Rand index of clusters using different clustering algorithms

	Profiles						Motifs				
	Kmeans	Fuzzy	SOM	Hier	RF		Kmeans	Fuzzy	SOM	Hier	RF
Kmeans	1	0.544	0.629	0.668	0.251	1	0.592	0.794	0.622	0.358	
Fuzzy	0.544	1	0.562	0.491	0.355	0.592	1	0.626	0.511	0.447	
SOM	0.629	0.562	1	0.49	0.287	0.794	0.626	1	0.591	0.33	
Hier	0.668	0.491	0.49	1	0.272	0.622	0.511	0.591	1	0.312	
RF	0.251	0.355	0.287	0.272	1	0.358	0.447	0.33	0.312	1	

The results from the kmeans algorithm using the motif variability data can be seen at Fig. 3. The cluster with 26 houses shows very little variability in the timing of their regular activities and can be assumed to be “creature of habit” households who may not respond well to an incentive to change behaviour. The 2 house and the 29 house clusters show lots of repeating activities and may be best to target for interventions as there are likely to be many activities that often repeat and that may be modifiable.

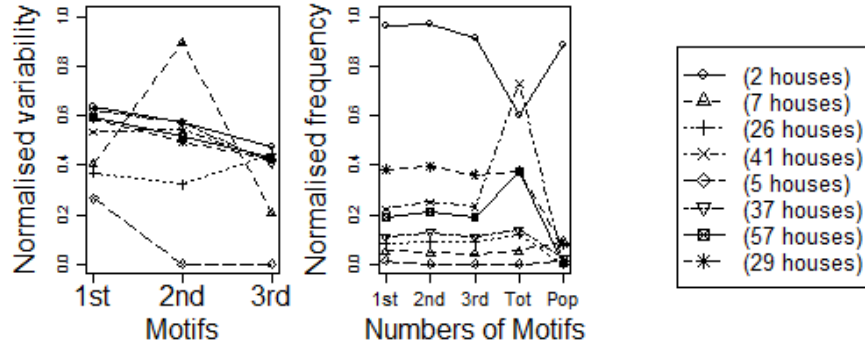


Fig. 3. kmeans clusters using motif variability

Examining the 29 house cluster in more detail, Fig. 4 show the motifs found for one of the houses and how the time of occurrence of the motif varies across the 4pm to 8pm period. In contrast, the motifs for one of the houses in the 26 house cluster are shown in Fig. 5 and the timings can be seen to be less variable.

As a comparison, the average load profiles for each of the households in the 29 house cluster are shown at Fig. 6. There is little similarity between the households and hence, using the load profile shapes as the basis, little likelihood of the households being clustered together. However, the variability in timing

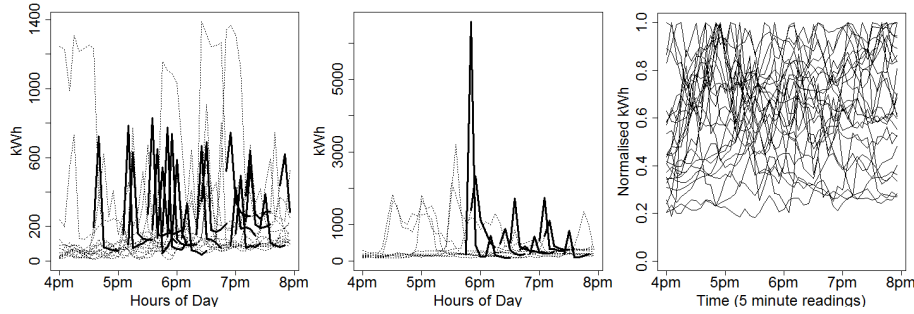


Fig. 4. Example house (high variability) **Fig. 5.** Example house (low variability) **Fig. 6.** Load profiles for high variability cluster

of the motifs can be used as a method for selecting appropriate households to target and allows groupings to be designated as high or low variability.

4 Significance and Impact

The ability to cost effectively partition domestic households into a few meaningful archetypes based on the household electricity usage is an important problem for the electricity industry. Identifying a few archetypal representations of households is essential for cost effective implementation of DSM techniques which itself is necessary to allow the electricity industry to meet the upcoming challenges. Producing more consistent and more descriptive archetypes than currently possible will allow the deployment of effective behaviour modification interventions.

Previous work does not incorporate any measure of the variability of regular behaviour when clustering households. The variability is an important characteristic as one of the major uses of the results is to target incentives for households to vary their behaviour to provide benefit to the electricity network.

The results presented show that the “variability in timing of motifs” approach produces more consistent clusters across different clustering algorithms compared to the consistency of clustering using just the daily load profiles.

The symbolisation technique is effective in detecting repeating patterns (motifs) that are approximately the same shape. Depending on the type of intervention planned for a subset of the households (for example, incentives to change overall electricity usage from day to night, or to influence short periods of usage during the peak period), different sizes of motifs may be used.

This work shows a novel approach to using electricity meter data to cluster households that enhances and complements the existing techniques based on the daily load profiles.

Acknowledgements

This work was possible thanks to RCUK Energy Programme and EPSRC grant references EP/I000496/1 and EP/G065802/1 and forms part of the Desimax project [27].

Thanks are due to Pavel Senin for providing R code implementing the SAX method.

References

1. DECC: Towards a Smarter Future, Government Response to the Consultation on Electricity and Gas Smart Metering. (2009)
2. River: Primer on demand-side management with an emphasis on price-responsive programs. prepared for The World Bank by Charles River Associates, Tech. Rep (2005)
3. Mooi, E., Sarstedt, M.: A concise guide to market research: The process, data, and methods using IBM SPSS statistics. Springer (2011)
4. Ramos, S., Figueiredo, V., Rodrigues, F., Pinheiro, R., Vale, Z.: Knowledge extraction from medium voltage load diagrams to support the definition of electrical tariffs. *International Journal of Engineering Intelligent Systems for Electrical Engineering and Communications* **15**(3) (2007) 143
5. Ellegård, K., Palm, J.: Visualizing energy consumption activities as a tool for making everyday life more sustainable. *Applied Energy* **88**(5) (2011) 1920–1926
6. Chicco, G.: Overview and performance assessment of the clustering methods for electrical load pattern grouping. *Energy* **Volume 42, Issue 1** (June 2012) 6880
7. Ramos, S., Duarte, J., Soares, J., Vale, Z., Duarte, F.: Typical load profiles in the smart grid context - A clustering methods comparison. In: *Power and Energy Society General Meeting, IEEE* (2012) 1–8
8. Lin, J., Keogh, E., Wei, L., Lonardi, S.: Experiencing SAX: a novel symbolic representation of time series. *Data Mining and Knowledge Discovery* **15**(2) (2007) 107–144
9. Shieh, J., Keogh, E.: i SAX: indexing and mining terabyte sized time series. In: *Proceeding of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining, ACM* (2008) 623–631
10. Mueen, A., Keogh, E., Zhu, Q., Cash, S., Westover, B.: Exact discovery of time series motifs. In: *Proc. of 2009 SIAM International Conference on Data Mining*. (2009) 1–12
11. Lines, J., Bagnall, A., Caiger-Smith, P., Anderson, S.: Classification of household devices by electricity usage profiles. *Intelligent Data Engineering and Automated Learning* (2011) 403–412
12. Craig, T., Polhill, J.G., Dent, I., Galan-Diaz, C., Heslop, S.: The North East Scotland Energy Monitoring Project: Exploring relationships between household occupants and energy usage. *Energy and Buildings* (2014)
13. Dent, I., Craig, T., Aickelin, U., Rodden, T.: A Method for Cleaning and Storing Electricity Meter Data for Flexible Analysis. In: *BeHave 2012, Helsinki*. (2012)
14. Elexon: The Electricity Trading Arrangements: A Beginners Guide. Technical report, Elexon (2012)
15. Jain, A., Dubes, R.: Algorithms for clustering data. Number 978-0130222787. Prentice Hall College Div (1988)

16. Bezdek, J.C.: Pattern recognition with fuzzy objective function algorithms. Kluwer Academic Publishers (1981)
17. Kohonen, T.: The self-organizing map. *Proceedings of the IEEE* **78**(9) (2002) 1464–1480
18. Everitt, B.S., Landau, S., Leese, M., Stahl, D.: Cluster analysis. Edward Arnold, London (2001)
19. Ward Jr, J.H.: Hierarchical grouping to optimize an objective function. *Journal of the American statistical association* **58**(301) (1963) 236–244
20. Breiman, L.: Random forests. *Machine learning* **45**(1) (2001) 5–32
21. Liaw, A., Wiener, M.: Classification and Regression by randomForest. *R News* **2**(3) (2002) 18–22
22. Electricity Association: Load profiles and their use in electricity settlement. UK-ERC (1997)
23. Figueiredo, V., Rodrigues, F., Vale, Z., Gouveia, J.: An electric energy consumer characterization framework based on data mining techniques. *Power Systems, IEEE Transactions on* **20**(2) (2005) 596–602
24. Chicco, G., Napoli, R., Postolache, P., Scutariu, M., Toader, C.: Customer characterization options for improving the tariff offer. *Power Systems, IEEE Transactions on* **18**(1) (2003) 381–387
25. Hubert, L., Arabie, P.: Comparing partitions. *Journal of classification* **2**(1) (1985) 193–218
26. Dent, I., Craig, T., Aickelin, U., Rodden, T.: Finding the creatures of habit; Clustering households based on their flexibility in using electricity. In: *Digital Futures*, Aberdeen, UK. (2012)
27. Kiprakis, A., Dent, I., Djokic, S., McLaughlin, S.: Multi-scale Dynamic Modeling to Maximize Demand Side Management. In: *IEEE Power and Energy Society Innovative Smart Grid Technologies Europe 2011*, Manchester, UK. (2011)